

Argonne Training Program on

EXTREME-SCALE COMPUTING



August 2 – August 14, 2015

Adaptive Linear Solvers and Eigensolvers

Jack Dongarra

University of Tennessee
Oak Ridge National Laboratory
University of Manchester

Dense Linear Algebra

.. Common Operations

$$Ax = b; \quad \min_x \|Ax - b\|; \quad Ax = \lambda x$$

.. A major source of large dense linear systems is problems involving the solution of boundary integral equations.

- The price one pays for replacing three dimensions with two is that what started as a sparse problem in $O(n^3)$ variables is replaced by a dense problem in $O(n^2)$.

.. Dense systems of linear equations are found in numerous other applications, including:

- airplane wing design;
- radar cross-section studies;
- flow around ships and other off-shore constructions;
- diffusion of solid bodies in a liquid;
- noise reduction; and
- diffusion of light through small particles.

Existing Math Software - Dense LA

DIRECT SOLVERS	License	Support	Type		Language			Mode		
			Real	Complex	F77/ F95	C	C++	Shared	Accel.	Dist
Chameleon	CeCILL-C	See authors	X	X		X		X	C	M
DPLASMA	BSD	yes	X	X		X		X	C	M
Eigen	Mozilla	yes	X	X			X	X		
Elemental	New BSD	yes	X	X			X			M
ELPA	LGPL	yes	X	X	F90	X		X		M
FLENS	BSD	yes	X	X			X	X		
hmat-oss	GPL	yes	X	X	X	X	X	X		
LAPACK	BSD	yes	X	X	X	X		X		
LAPACK95	BSD	yes	X	X	X			X		
libflame	New BSD	yes	X	X	X	X		X		
MAGMA	BSD	yes	X	X	X	X		X	C/O/X	
NAPACK	BSD	yes	X		X			X		
PLAPACK	LGPL	yes	X	X	X	X				M
PLASMA	BSD	yes	X	X	X	X		X		
rejtrix	by-nc-sa	yes	X				X	X		
ScaLAPACK	BSD	yes	X	X	X	X				M/P
Trilinos/Pliris	BSD	yes	X	X		X	X			M
ViennaCL	MIT	yes	X				X	X	C/O/X	

<http://www.netlib.org/utk/people/JackDongarra/la-sw.html>

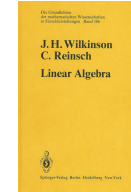
.. LINPACK, EISPACK, LAPACK, ScaLAPACK

8/10/15 ➤ PLASMA, MAGMA

DLA Solvers

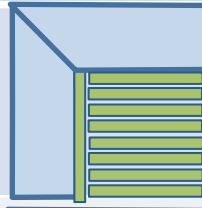
- We are interested in developing Dense Linear Algebra Solvers
- Retool LAPACK and ScaLAPACK for multicore and hybrid architectures

40 Years Evolving SW and Alg Tracking Hardware Developments



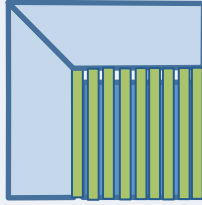
Software/Algorithms follow hardware evolution in time

EISPACK (70's)
(Translation of Algol)



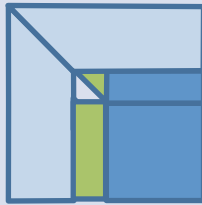
Rely on
- Fortran, but row oriented

LINPACK (80's)
(Vector operations)



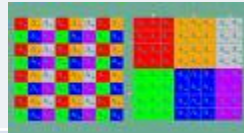
Rely on
- Level-1 BLAS operations
- Column oriented

LAPACK (90's)
(Blocking, cache friendly)



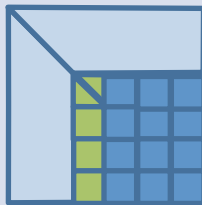
Rely on
- Level-3 BLAS operations

ScaLAPACK (00's)
(Distributed Memory)



Rely on
- PBLAS Mess Passing

PLASMA (10's)
New Algorithms
(many-core friendly)



Rely on
- DAG/scheduler
- block data layout
- some extra kernels

What do you mean by performance?

♦ What is a flop/s?

- flop/s is a rate of execution, some number of floating point operations per second.
 - » Whenever this term is used it will refer to 64 bit floating point operations and the operations will be either addition or multiplication.

♦ What is the theoretical peak performance?

- The theoretical peak is based not on an actual performance from a benchmark run, but on a paper computation to determine the theoretical peak rate of execution of floating point operations for the machine.
- The theoretical peak performance is determined by counting the number of floating-point additions and multiplications (in full precision) that can be completed during a period of time, usually the cycle time of the machine.
- For example, an Intel Xeon Haswell dual core at 2.3 GHz can complete 16 floating point operations per cycle or a theoretical peak performance of 36.8 GFlop/s per core or 73.6 Gflop/s for the socket.

Peak Performance - Per Core

$$\text{FLOPS} = \text{cores} \times \text{clock} \times \frac{\text{FLOPs}}{\text{cycle}}$$

Floating point operations per cycle per core

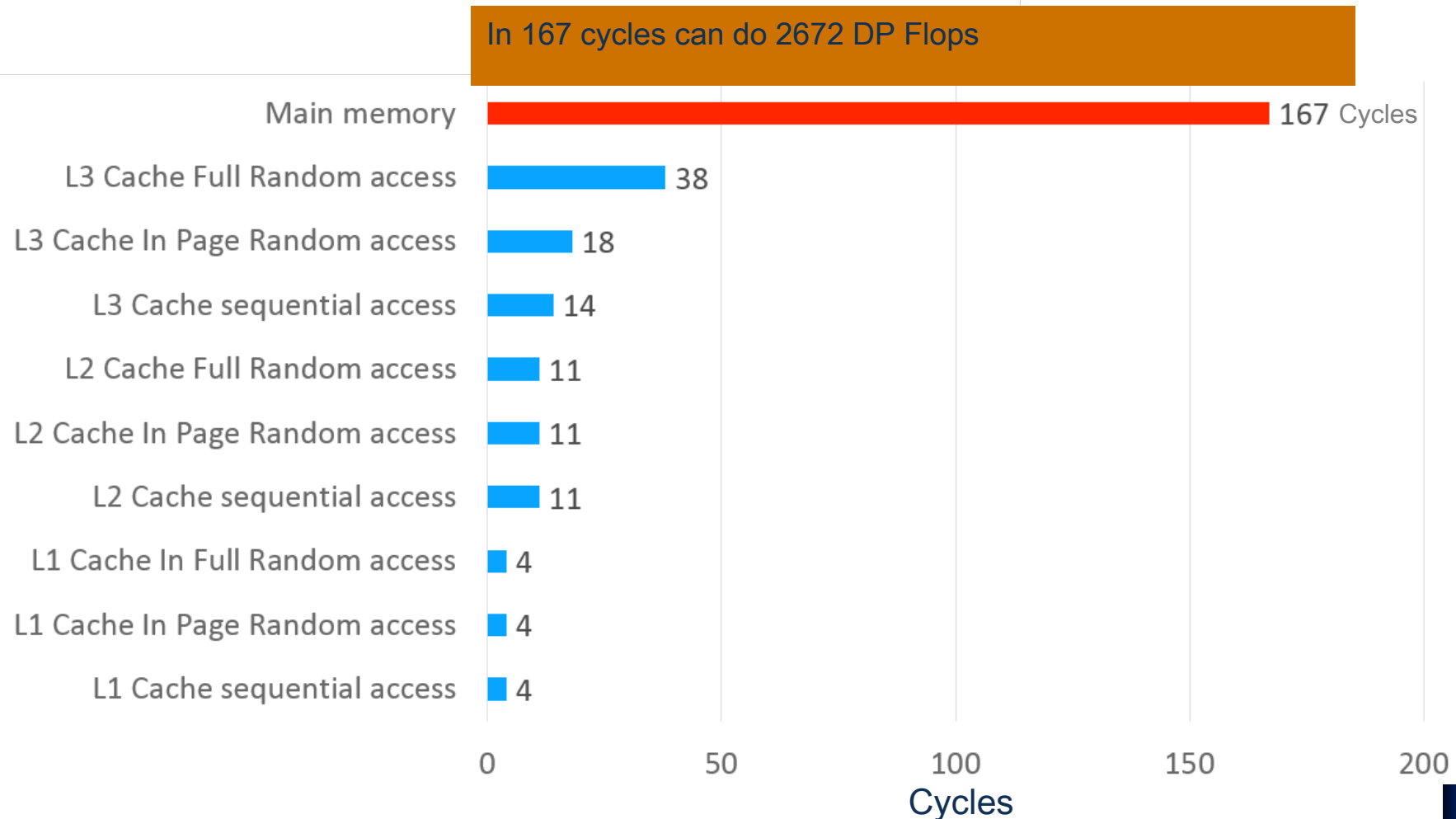
- + Most of the recent computers have FMA (Fused multiple add): (i.e. $x \leftarrow x + y * z$ in one cycle)
- + Intel Xeon earlier models and AMD Opteron have SSE2
 - + 2 flops/cycle DP & 4 flops/cycle SP
- + Intel Xeon Nehalem ('09) & Westmere ('10) have SSE4
 - + 4 flops/cycle DP & 8 flops/cycle SP
- + Intel Xeon Sandy Bridge('11) & Ivy Bridge ('12) have AVX & AVX2
 - + 8 flops/cycle DP & 16 flops/cycle SP
- ➔ + Intel Xeon Haswell ('13) & (Broadwell ('14)) AVX2
 - + 16 flops/cycle DP & 32 flops/cycle SP
 - + Xeon Phi (per core) is at 16 flops/cycle DP & 32 flops/cycle SP
- + Intel Xeon Skylake ('15, Q3)
 - + 32 flops/cycle DL & 64 flops/cycle SP



We
are
here

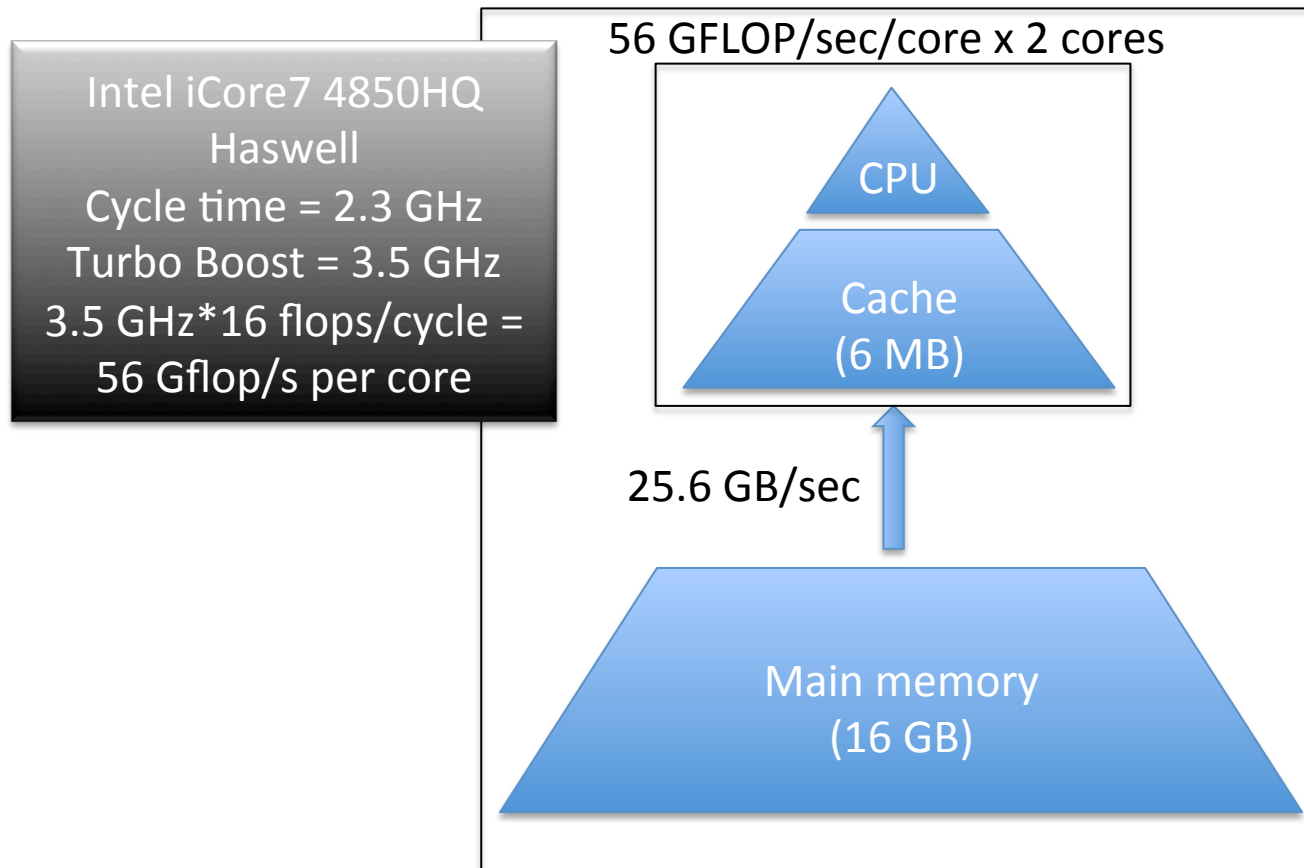
87 GFLOPS (DP-FP, peak)	185 GFLOPS (DP-FP, peak)	~225 GFLOPS (DP-FP, peak)	~500 GFLOPS (DP-FP, peak)	1tb GFLOPS (DP-FP, peak)	1tb GFLOPS (DP-FP, peak)
Westmere	Sandy Bridge	Ivy Bridge	Haswell	Broadwell	Skylake ...
32nm SSE4.2 DDR3 PCIe2	32nm AVX DDR3 PCIe3	22nm	22nm AVX2 DDR4 PCIe3	14nm	14nm AVX2 DDR4 PCIe4

CPU Access Latencies in Clock Cycles



Memory transfer

- One level of memory model on my laptop:

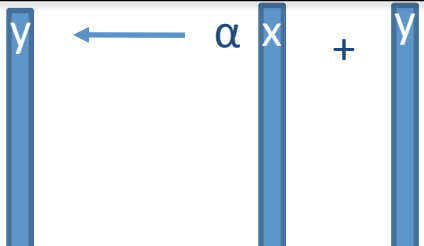


(Omitting latency here.)

The model IS simplified (see next slide) but it provides an upper bound on performance as well. I.e., we will never go faster than what the model predicts.
(And, of course, we can go slower ...)

8/10/15

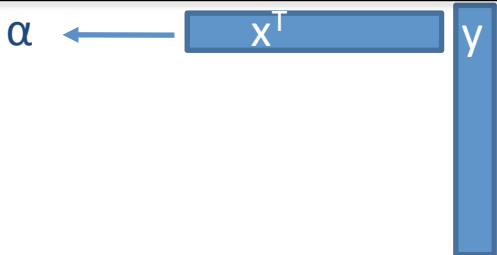
FMA: fused multiply-add

AXPY:  $\alpha x + y$

```
for ( j = 0; j < n; j++)  
    y[i] += a * x[i];
```

(without increment)

n MUL
n ADD
2n FLOP
n FMA

DOT:  $\alpha \leftarrow x^T y$

```
alpha = 0e+00;  
for ( j = 0; j < n; j++)  
    alpha += x[i] * y[i];
```

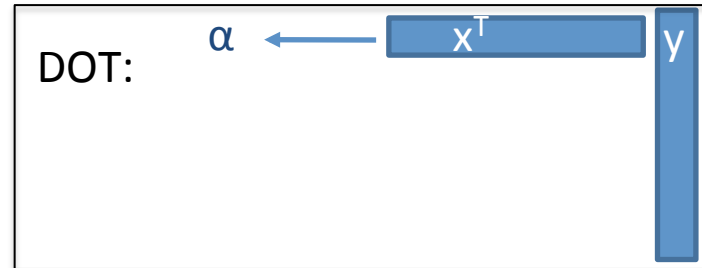
(without increment)

n MUL
n ADD
2n FLOP
n FMA

Note: It is reasonable to expect the one loop codes shown here to perform as well as their Level 1 BLAS counterpart (on multicore with an OpenMP pragma for example).

The true gain these days with using the BLAS is (1) Level 3 BLAS, and (2) portability.

- Take two double precision vectors x and y of size $n=375,000$.



- Data size:
 - $(375,000 \text{ double}) * (8 \text{ Bytes / double}) = 3 \text{ MBytes}$ per vector
 - (Two vectors fit in cache (6 MBytes). OK.)

- Time to move the vectors from memory to cache:
 - $(6 \text{ MBytes}) / (25.6 \text{ GBytes/sec}) = \mathbf{0.23 \text{ ms}}$
- Time to perform computation of DOT:
 - $(2n \text{ flop}) / (56 \text{ Gflop/sec}) = \mathbf{0.01 \text{ ms}}$

Vector Operations

$$\begin{aligned}\text{total_time} &\geq \max (\text{time_comm} , \text{time_comp}) \\ &= \max (0.23\text{ms} , 0.01\text{ms}) = 0.23\text{ms}\end{aligned}$$

$$\text{Performance} = (2 \times 375,000 \text{ flops}) / .23\text{ms} = 3.2 \text{ Gflop/s}$$

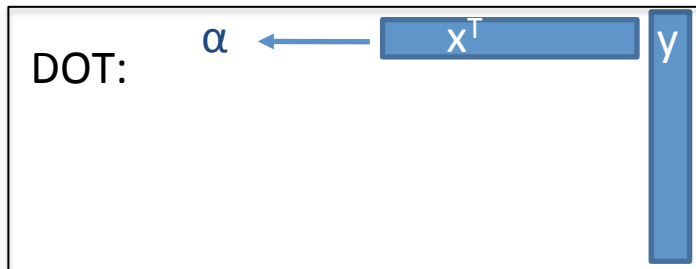
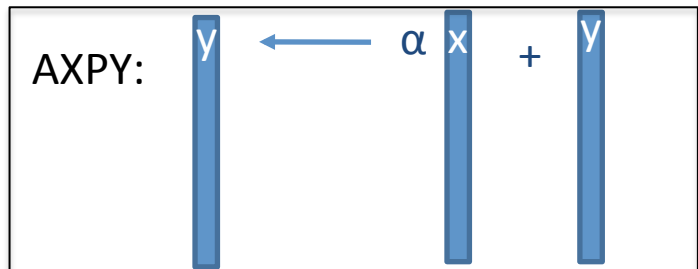
Performance for DOT ≤ 3.2 Gflop/s

Peak is 56 Gflop/s

We say that the operation is communication bounded. No reuse of data.

Level 1, 2 and 3 BLAS

Level 1 BLAS Matrix-Vector operations



2n FLOP

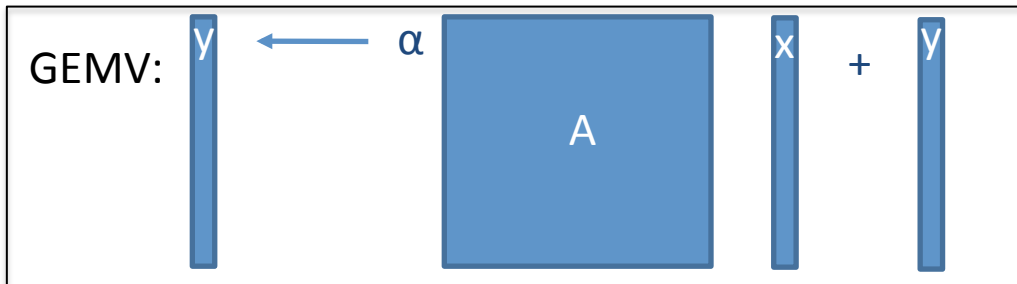
2n memory reference

AXPY: 2n READ, n WRITE

DOT: 2n READ

RATIO: 1

Level 2 BLAS Matrix-Vector operations

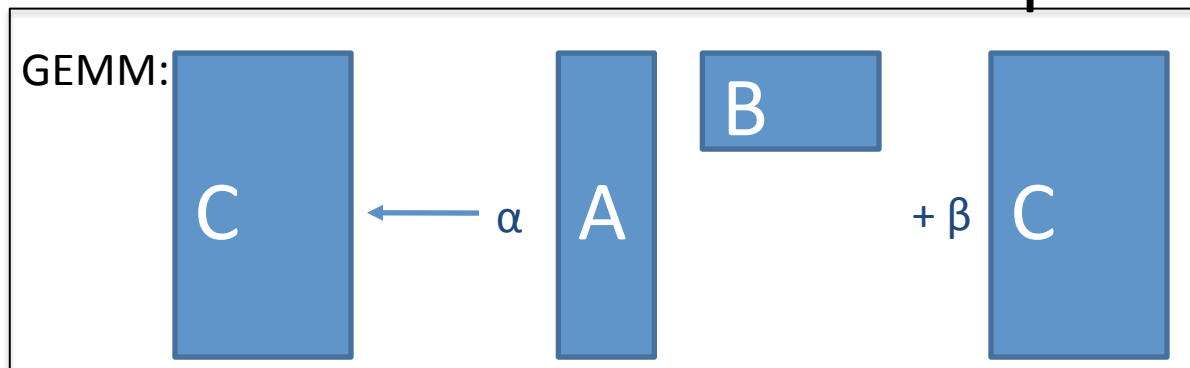


2n² FLOP

n² memory references

RATIO: 2

Level 3 BLAS Matrix-Matrix operations



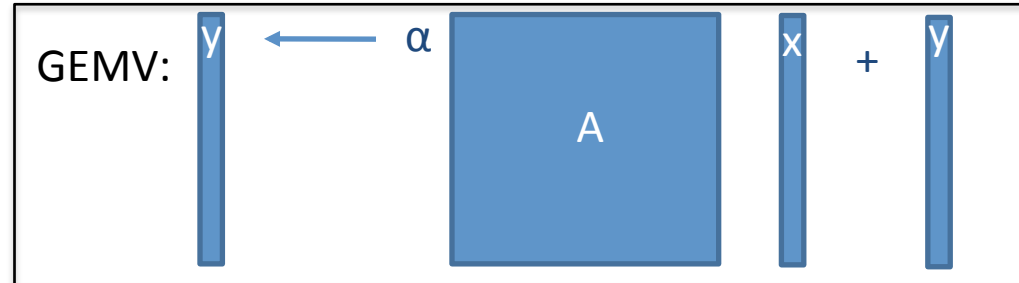
2n³ FLOP

3n² memory references

3n² READ, n² WRITE

RATIO: 2/3 n

- Double precision matrix A and vectors x and y of size $n=860$.



- Data size:

$$- (860^2 + 2*860 \text{ double}) * (8 \text{ Bytes} / \text{double}) \sim 6 \text{ MBytes}$$

Matrix and two vectors fit in cache (6 MBytes).

- Time to move the data from memory to cache:

$$- (6 \text{ MBytes}) / (25.6 \text{ GBytes/sec}) = \mathbf{0.23 \text{ ms}}$$

- Time to perform computation of DOT:

$$- (2n^2 \text{ flop}) / (56 \text{ Gflop/sec}) = \mathbf{0.26 \text{ ms}}$$

Matrix - Vector Operations

$$\begin{aligned}\text{total_time} &\geq \max (\text{time_comm} , \text{time_comp}) \\ &= \max (0.23\text{ms} , 0.26\text{ms}) = 0.26\text{ms}\end{aligned}$$

$$\text{Performance} = (2 \times 860^2 \text{ flops}) / .26\text{ms} = 5.7 \text{ Gflop/s}$$

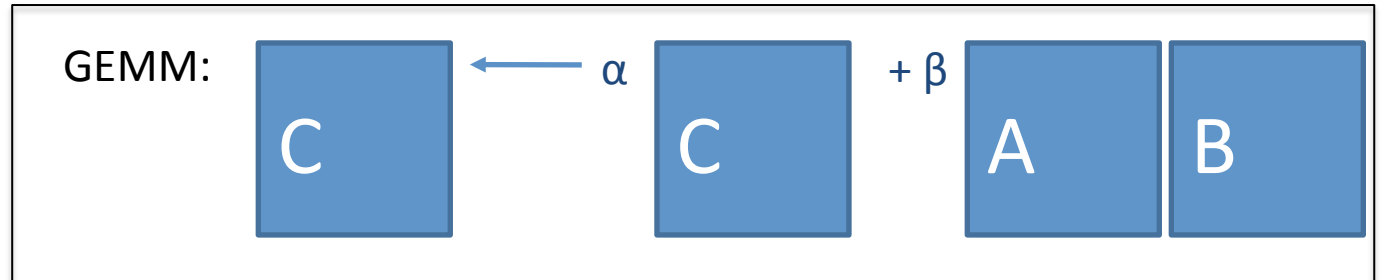
Performance for GEMV ≤ 5.7 Gflop/s

Performance for DOT ≤ 3.2 Gflop/s

Peak is 56 Gflop/s

We say that the operation is communication bounded. Very little reuse of data.

- Take two double precision vectors x and y of size $n=500$.



- Data size:
 - $(500^2 \text{ double}) * (8 \text{ Bytes / double}) = 2 \text{ MBytes per matrix}$
 - (Three matrices fit in cache (6 MBytes). OK.)
- Time to move the matrices in cache:
 - $(6 \text{ MBytes}) / (25.6 \text{ GBytes/sec}) = \mathbf{0.23 \text{ ms}}$
- Time to perform computation in GEMM:
 - $(2n^3 \text{ flop}) / (56 \text{ Gflop/sec}) = \mathbf{4.46 \text{ ms}}$

Matrix Matrix Operations

$$\begin{aligned}\text{total_time} &\geq \max(\text{time_comm}, \text{time_comp}) \\ &= \max(0.23\text{ms}, 4.46\text{ms}) = 4.46\text{ms}\end{aligned}$$

For this example, communication time is less than 6% of the computation time.

$$\text{Performance} = (2 \times 500^3 \text{ flops}) / 4.69\text{ms} = 53.3 \text{ Gflop/s}$$

There is a lots of data reuse in a GEMM; $2/3n$ per data element. Has good temporal locality.

If we assume $\text{total_time} \approx \text{time_comm} + \text{time_comp}$, we get

Performance for GEMM ≈ 53.3 Gflop/sec

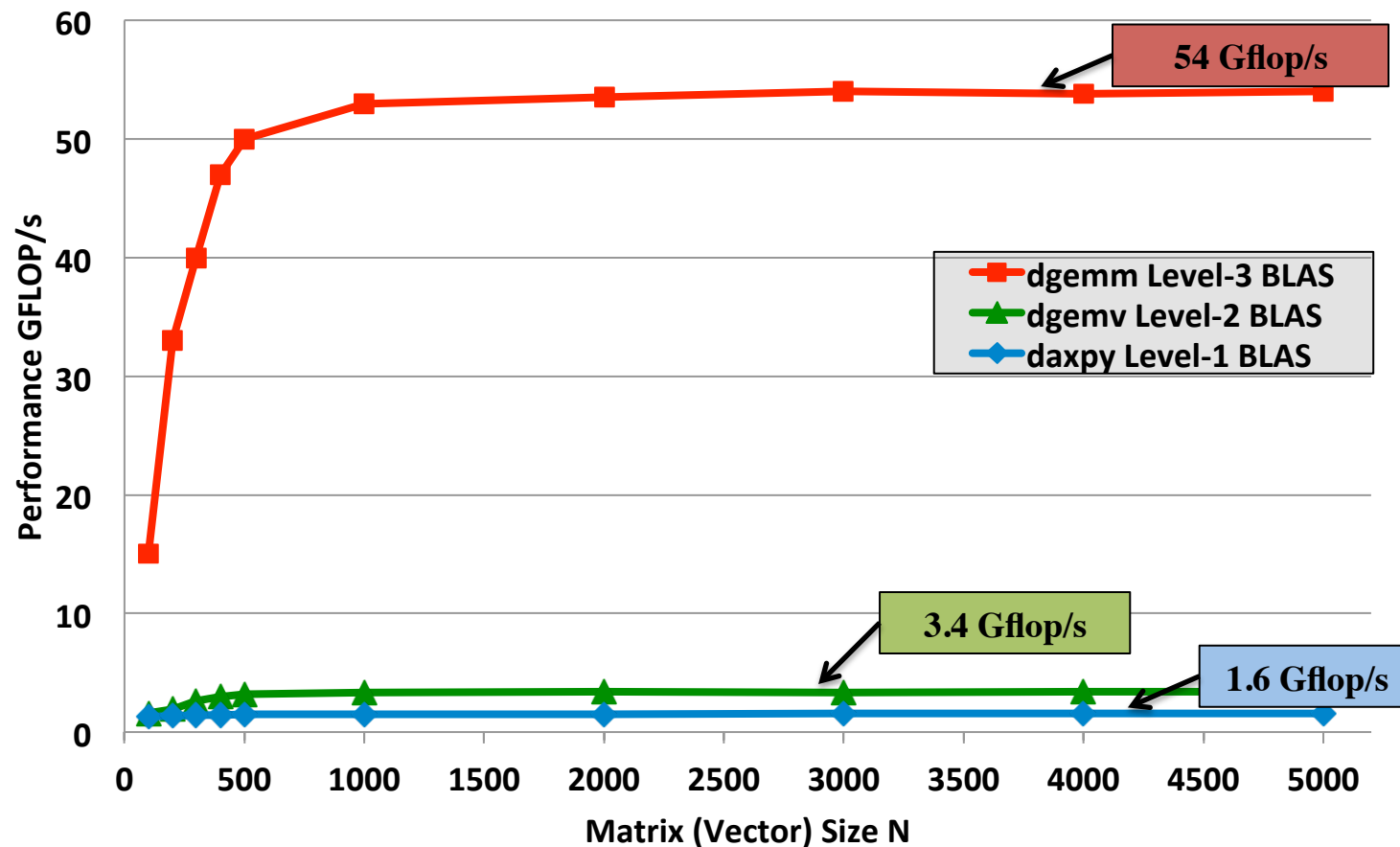
Performance for DOT ≤ 3.2 Gflop/s

Performance for GEMV ≤ 5.7 Gflop/s

(Out of 56 Gflop/sec possible, so that would be 95% peak performance efficiency.)

Level 1, 2 and 3 BLAS

1 core Intel Haswell i7-4850HQ, 2.3 GHz (Turbo Boost at 3.5 GHz);
Peak = 56 Gflop/s



1 core Intel Haswell i7-4850HQ, 2.3 GHz, Memory: DDR3L-1600MHz
6 MB shared L3 cache, and each core has a private 256 KB L2 and 64 KB L1.
The theoretical peak per core double precision is 56 Gflop/s per core.
Compiled with gcc and using VecLib

Issues

- Reuse based on matrices that fit into cache.
- What if you have matrices bigger than cache?

Issues

- Reuse based on matrices that fit into cache.
- What if you have matrices bigger than cache?
- Break matrices into blocks or tiles that will fit.

The diagram illustrates a sequence of matrix operations using 3x3 grids of tiles. Each tile is a blue square containing a label. The first grid on the left represents matrix C, with tiles labeled C_{11} through C_{33} . An arrow points from this grid to the right, where it is followed by the scalar β . This is followed by another 3x3 grid of tiles labeled C_{11} through C_{33} . To the right of this grid is the scalar $+$ α , followed by a 3x3 grid of tiles labeled A_{11} through A_{33} . To the right of this grid is the scalar $*$, followed by a final 3x3 grid of tiles labeled B_{11} through B_{33} .

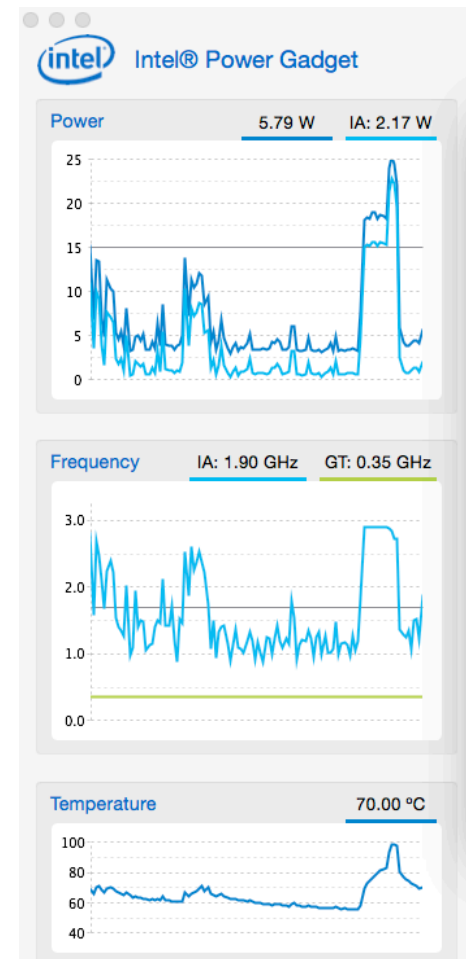
By the way

Performance for your laptop

- If you are interested in running the Linpack Benchmark on your system see:

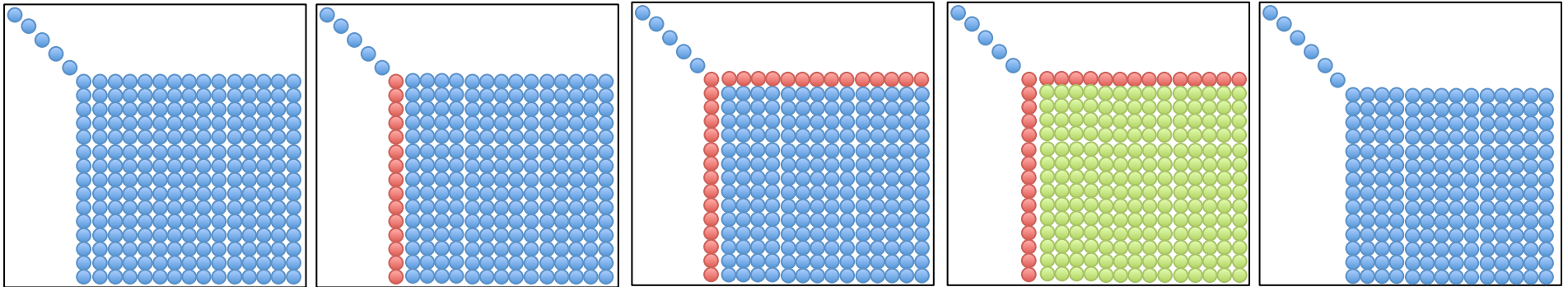
<https://software.intel.com/en-us/node/157667?wapkw=mkl+linpack>

- Also Intel has a power meter, see:
<https://software.intel.com/en-us/articles/intel-power-gadget-20>



The Standard LU Factorization LINPACK

1970's HPC of the Day: Vector Architecture



Factor column
with Level 1
BLAS

Divide by
Pivot
row

Schur
complement
update
(Rank 1 update)

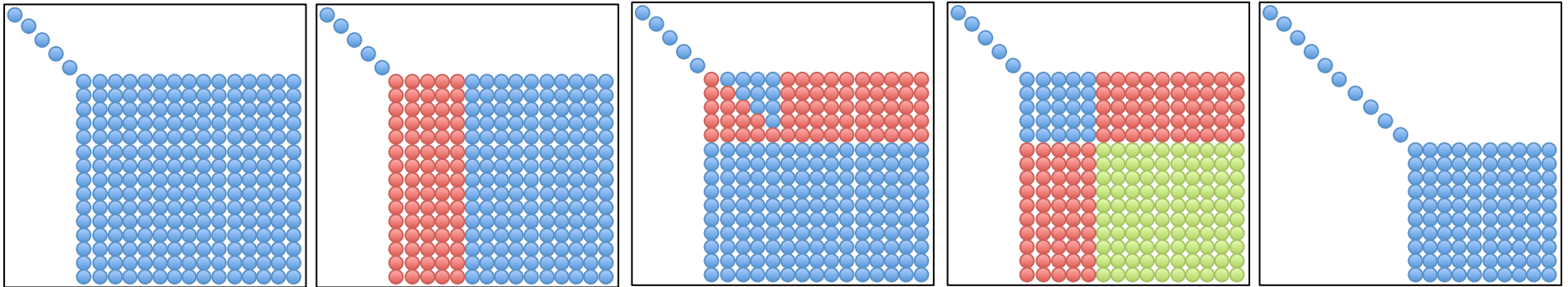
Next Step

Main points

- Factorization column (zero) mostly sequential due to memory bottleneck
- Level 1 BLAS
- Divide pivot row has little parallelism
- Rank -1 Schur complement update is the only easy parallelize task
- Partial pivoting complicates things even further
- Bulk synchronous parallelism (fork-join)
 - Load imbalance
 - Non-trivial Amdahl fraction in the panel
 - Potential workaround (look-ahead) has complicated implementation

The Standard LU Factorization LAPACK

1980's HPC of the Day: Cache Based SMP



Factor panel
with Level 1,2
BLAS

Triangular
update

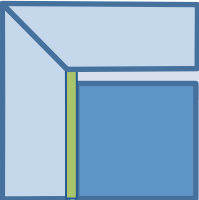

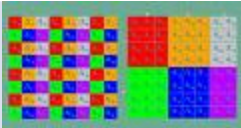
Schur
complement
update

Next Step

Main points

- Panel factorization mostly sequential due to memory bottleneck
- Triangular solve has little parallelism
- Schur complement update is the only easy parallelize task
- Partial pivoting complicates things even further
- Bulk synchronous parallelism (fork-join)
 - Load imbalance
 - Non-trivial Amdahl fraction in the panel
 - Potential workaround (look-ahead) has complicated implementation

Last Generations of DLA Software

Software/Algorithms follow hardware evolution in time		
LINPACK (70's) (Vector operations)		Rely on - Level-1 BLAS operations
LAPACK (80's) (Blocking, cache friendly)		Rely on - Level-3 BLAS operations
ScaLAPACK (90's) (Distributed Memory)		Rely on - PBLAS Mess Passing

2D Block Cyclic Layout

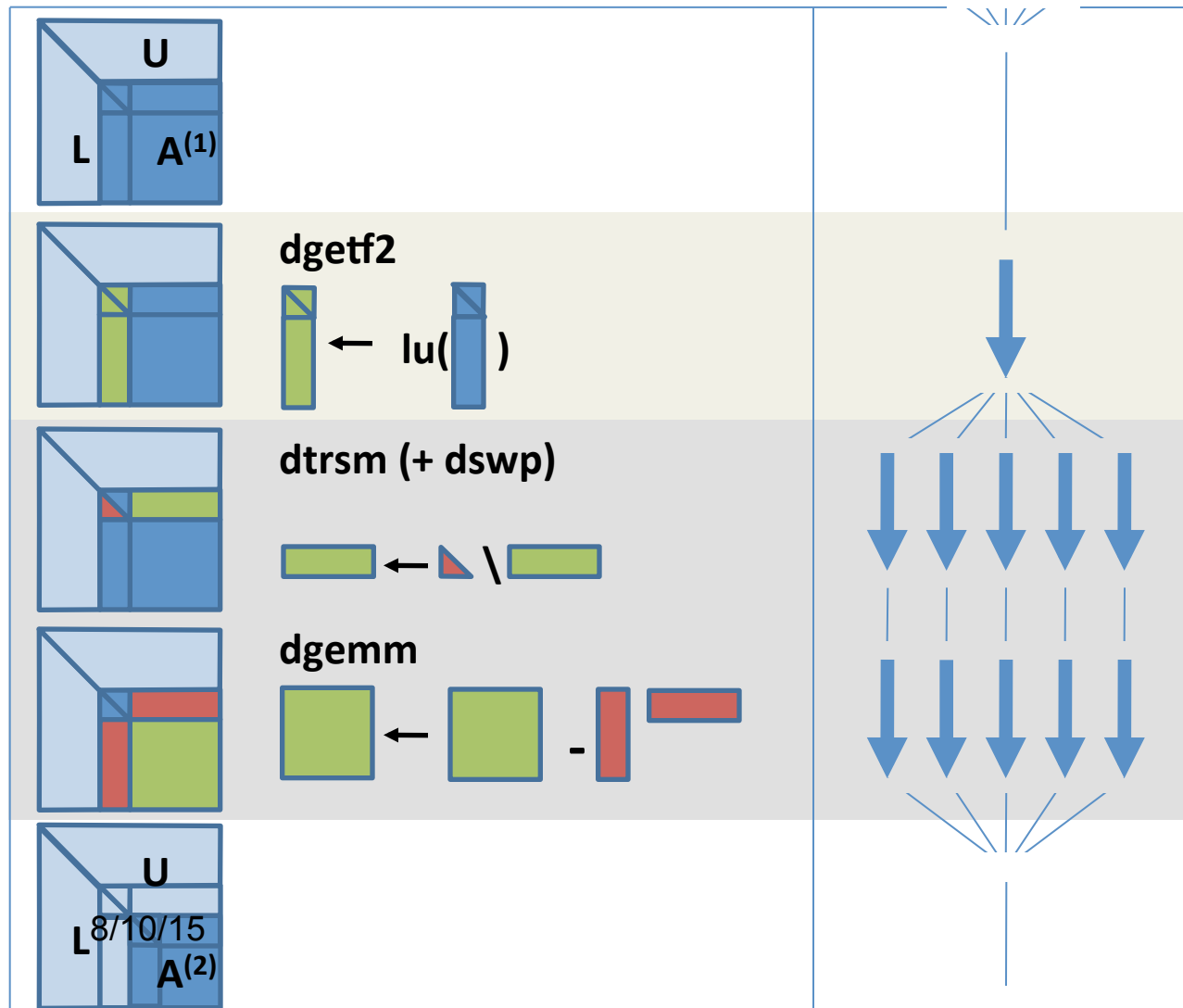
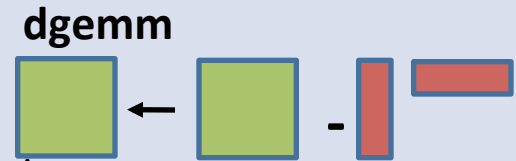
Matrix point of view									Processor point of view								
0	2	4	0	2	4	0	2	4	0	0	0	2	2	2	4	4	4
1	3	5	1	3	5	1	3	5	0	0	0	2	2	2	4	4	4
0	2	4	0	2	4	0	2	4	0	0	0	2	2	2	4	4	4
1	3	5	1	3	5	1	3	5	0	0	0	2	2	2	4	4	4
0	2	4	0	2	4	0	2	4	0	0	0	2	2	2	4	4	4
1	3	5	1	3	5	1	3	5	1	1	1	3	3	3	5	5	5
0	2	4	0	2	4	0	2	4	1	1	1	3	3	3	5	5	5
1	3	5	1	3	5	1	3	5	1	1	1	3	3	3	5	5	5
0	2	4	0	2	4	0	2	4	1	1	1	3	3	3	5	5	5

Parallelization of LU and QR.

IC

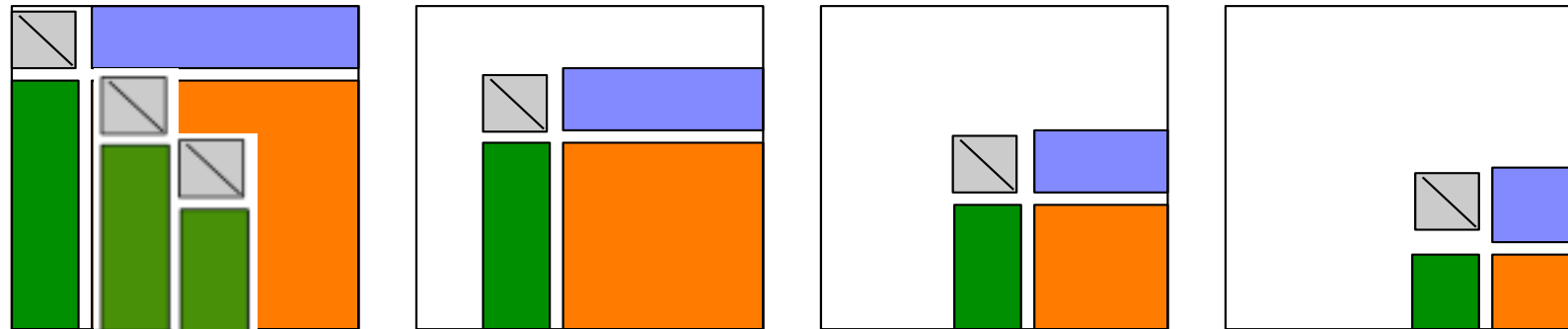
Parallelize the update:

- Easy and done in any reasonable software.
- This is the $2/3n^3$ term in the FLOPs count.
- Can be done efficiently with LAPACK+multithreaded BLAS

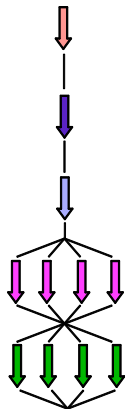
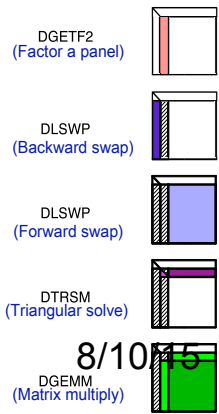
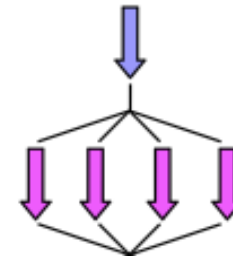
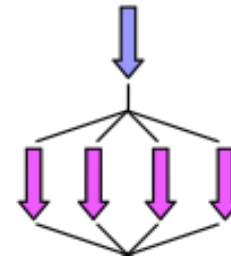
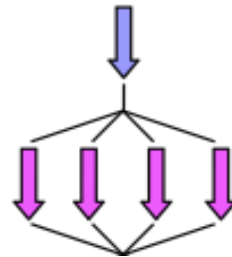
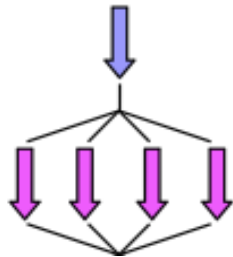


Fork - Join parallelism
Bulk Sync Processing

Synchronization (in LAPACK LU)



Step 1 → Step 2 → Step 3 → Step 4 . . .



LAPACK

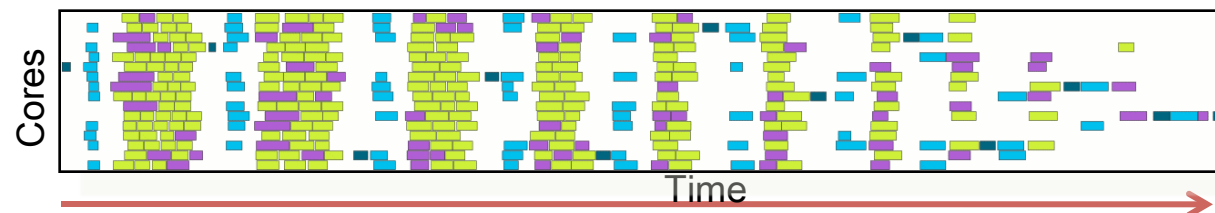
LAPACK

LAPACK

BLAS

BLAS

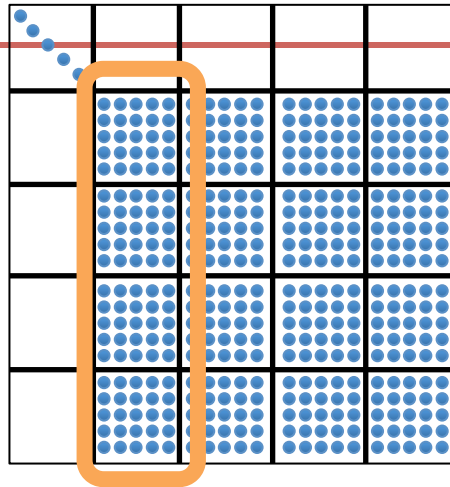
- fork join
- bulk synchronous processing



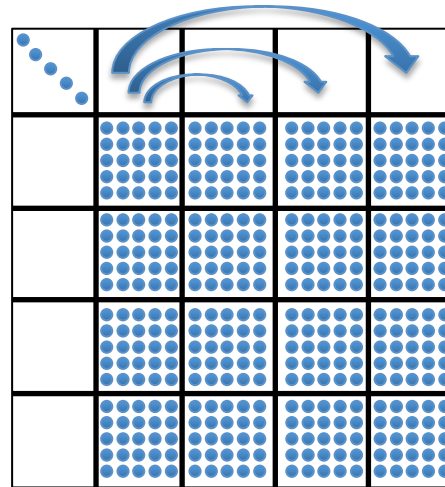
PLASMA LU Factorization

Dataflow Driven

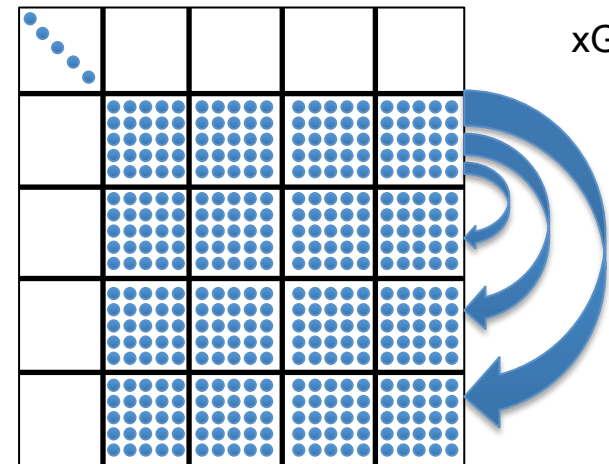
Numerical program generates tasks and run time system executes tasks respecting data dependences.



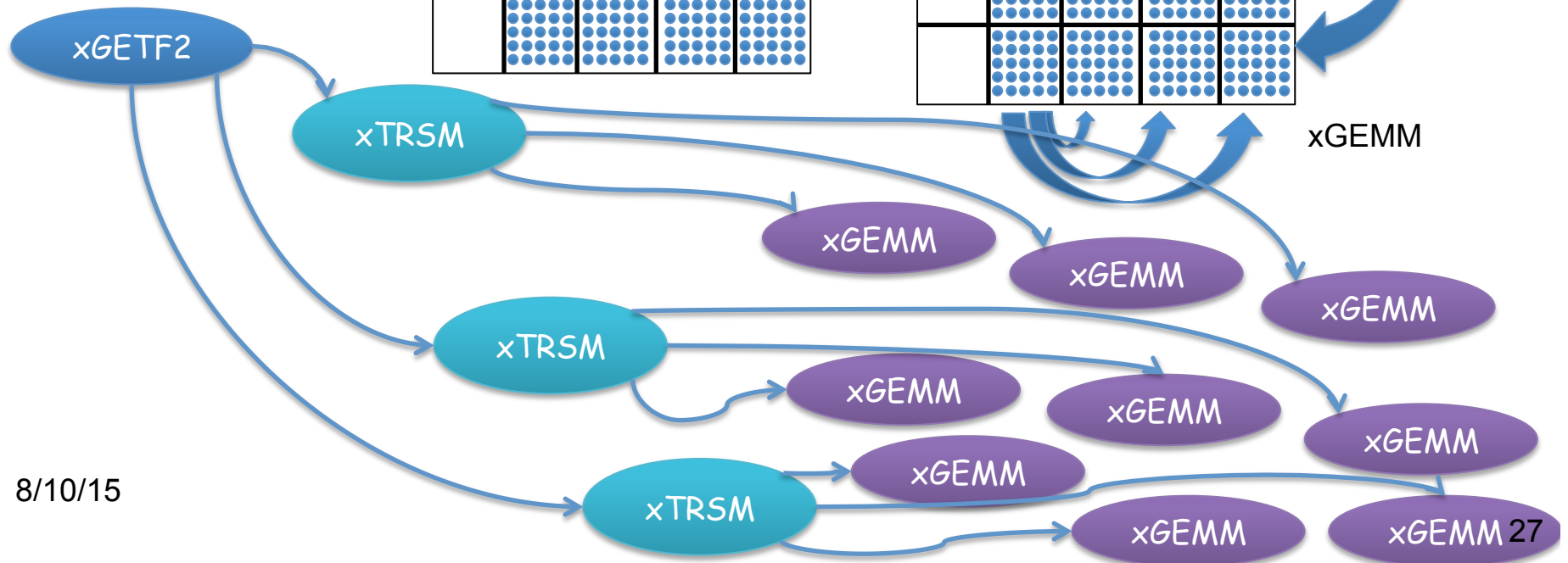
xTRSM



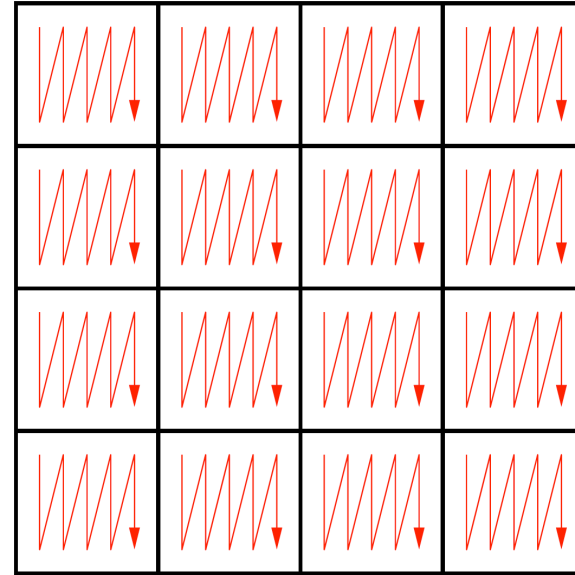
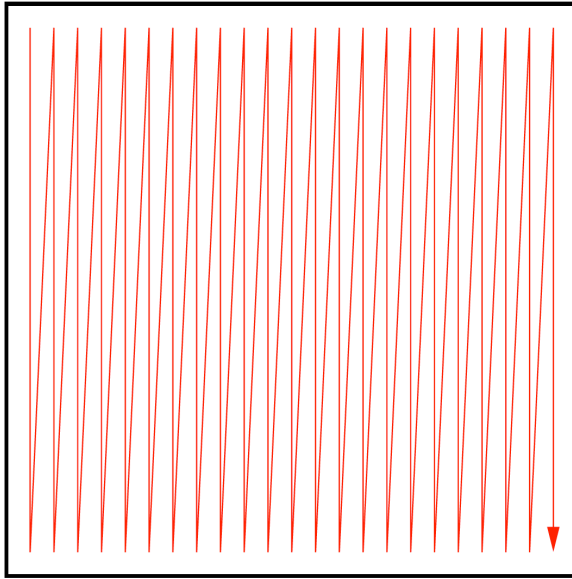
xGEMM



xGEMM



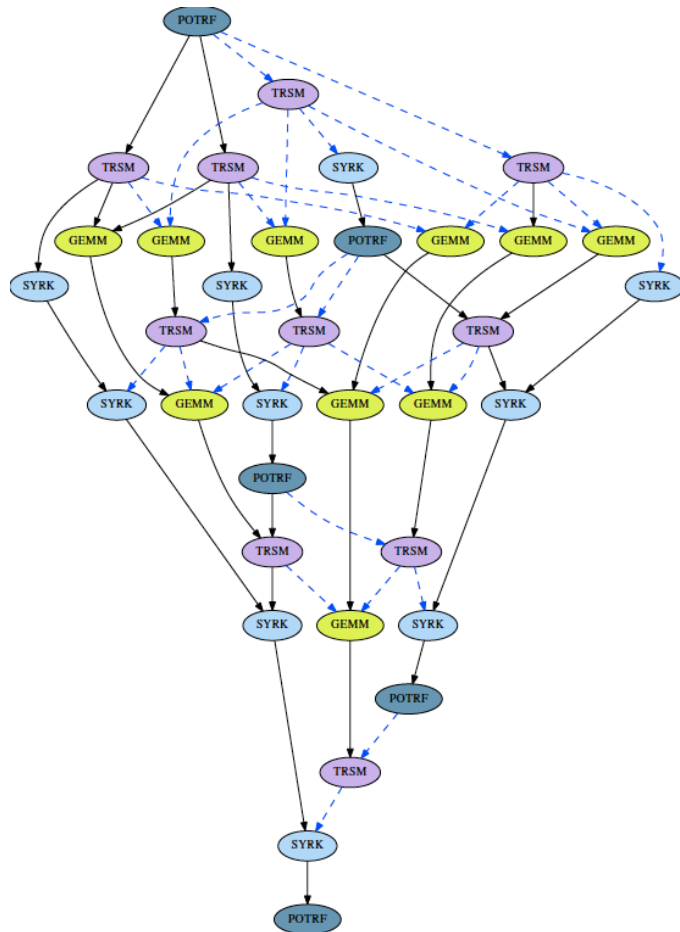
Data Layout is Critical



- Tile data layout where each data tile is contiguous in memory
- Decomposed into several fine-grained tasks, which better fit the memory of the small core caches

```
FOR k = 0..TILES-1
  A[k][k] ← DPOTRF(A[k][k])
  FOR m = k+1..TILES-1
    A[m][k] ← DTRSM(A[k][k], A[m][k])
  FOR m = k+1..TILES-1
    A[m][m] ← DSYRK(A[m][k], A[m][m])
    FOR n = k+1..m-1
      A[m][n] ← DGEMM(A[m][k], A[n][k], A[m][n])
```

definition – pseudocode



A runtime environment for the dynamic execution of precedence-constraint tasks (DAGs) in a multicore machine

- Translation
- If you have a serial program that consists of computational kernels (tasks) that are related by data dependencies, QUARK can help you execute that program (relatively efficiently and easily) in parallel on a multicore machine

The Purpose of a QUARK Runtime

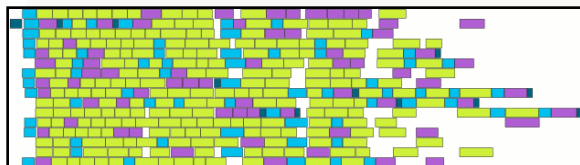
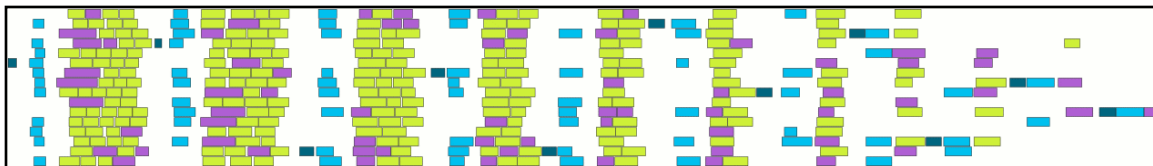
Objectives

- High utilization of each core
- Scaling to large number of cores
- Synchronization reducing algorithms

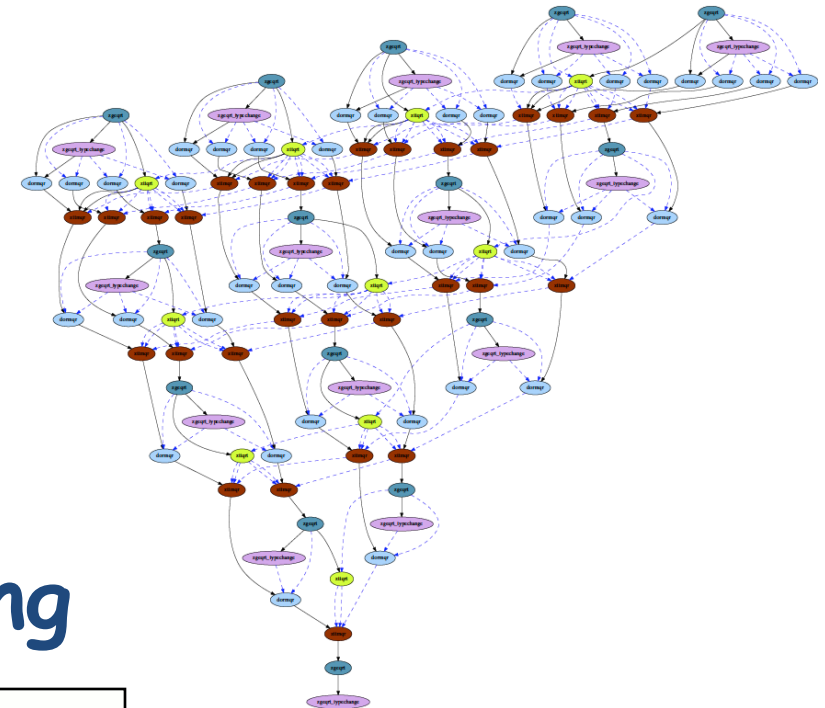
Methodology

- Dynamic DAG scheduling (QUARK)
- Explicit parallelism
- Implicit communication
- Fine granularity / block data layout

Arbitrary DAG with dynamic scheduling



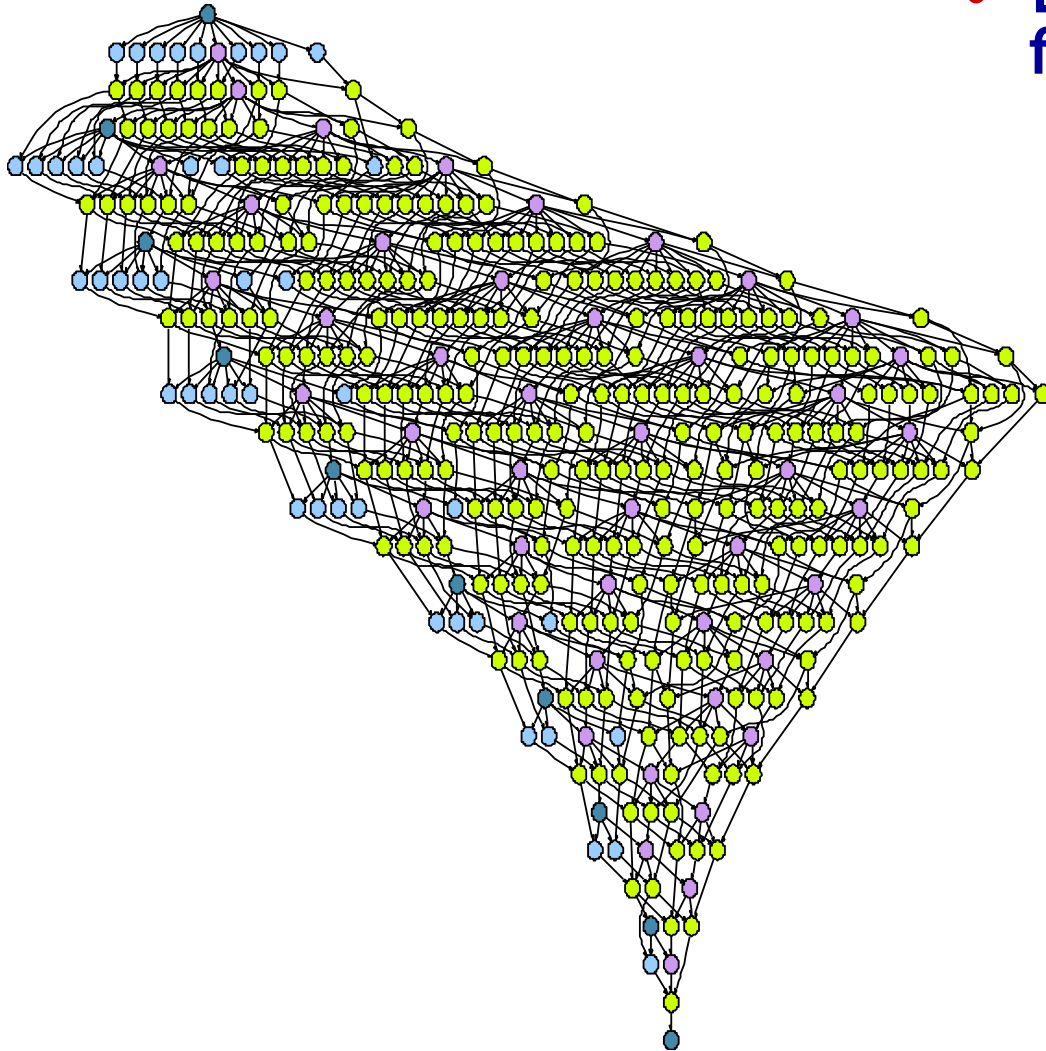
DAG scheduled parallelism



Fork-join parallelism
Notice the synchronization penalty in the presence of heterogeneity.

PLASMA Local Scheduling

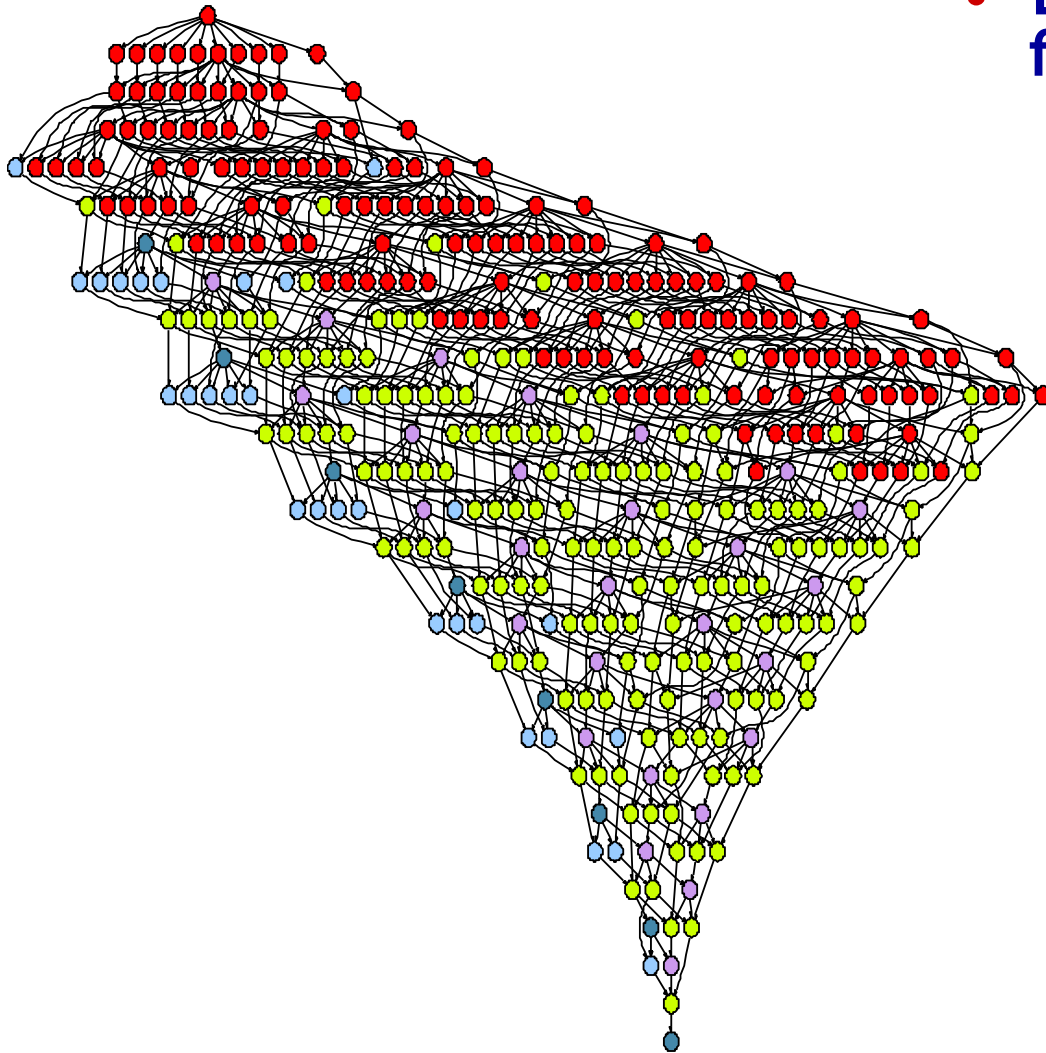
Dynamic Scheduling: Sliding Window



- **DAGs get very big, very fast**
 - **So windows of active tasks are used; this means no global critical path**
 - **Matrix of $NB \times NB$ tiles; NB^3 operation**
 - **$NB=100$ gives 1 million tasks**

PLASMA Local Scheduling

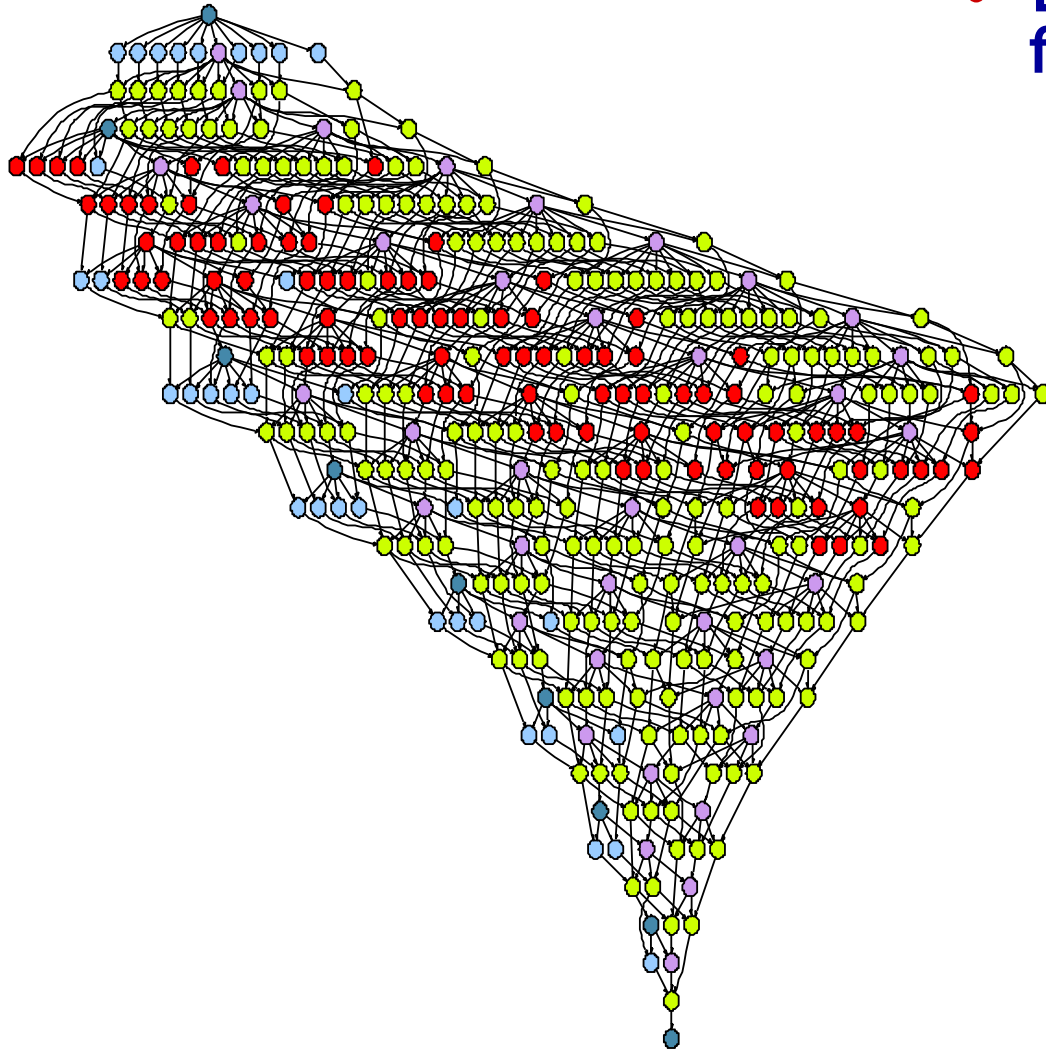
Dynamic Scheduling: Sliding Window



- DAGs get very big, very fast
 - So windows of active tasks are used; this means no global critical path
 - Matrix of $NB \times NB$ tiles; NB^3 operation
 - $NB=100$ gives 1 million tasks

PLASMA Local Scheduling

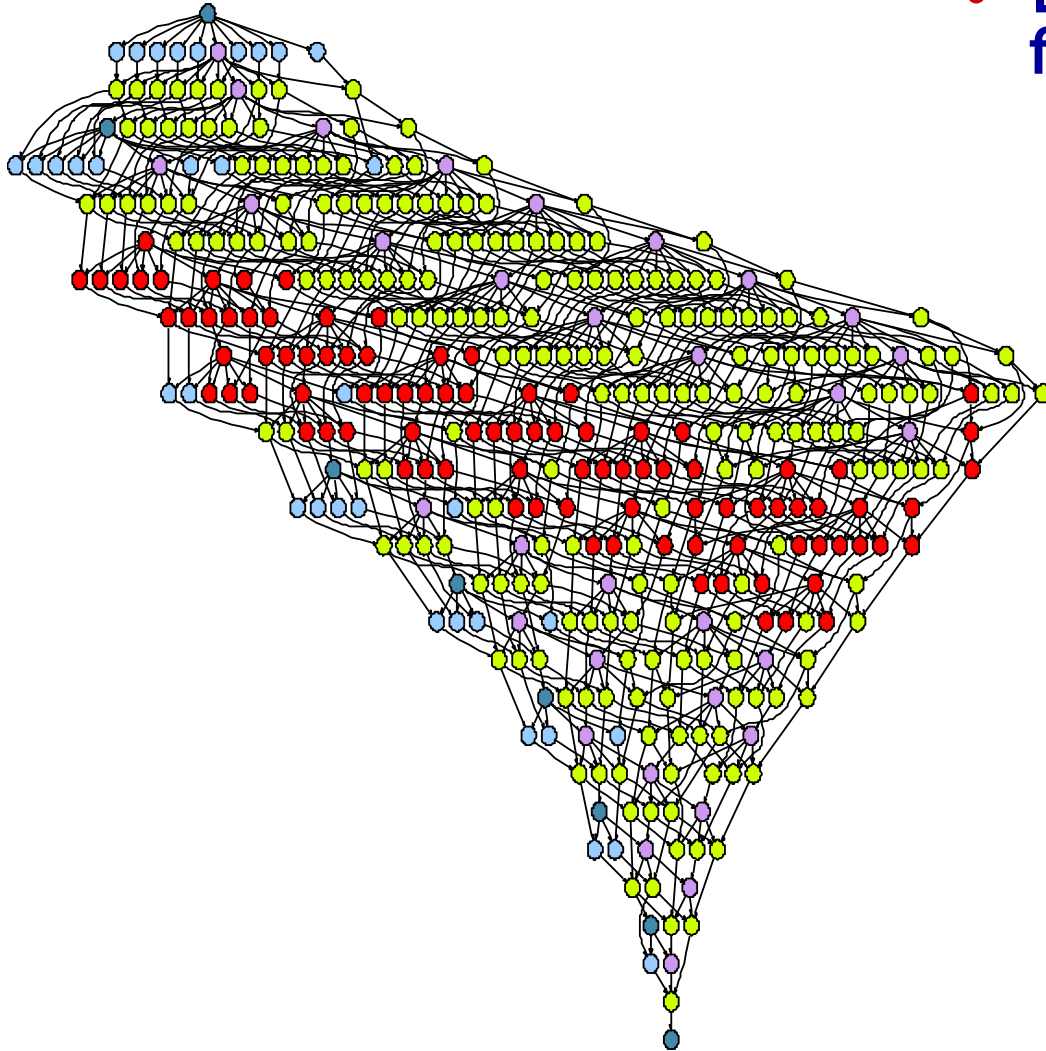
Dynamic Scheduling: Sliding Window



- DAGs get very big, very fast
 - So windows of active tasks are used; this means no global critical path
 - Matrix of $NB \times NB$ tiles; NB^3 operation
 - $NB=100$ gives 1 million tasks

PLASMA Local Scheduling

Dynamic Scheduling: Sliding Window



- **DAGs get very big, very fast**
 - So windows of active tasks are used; this means no global critical path
 - Matrix of $NB \times NB$ tiles; NB^3 operation
 - $NB=100$ gives 1 million tasks

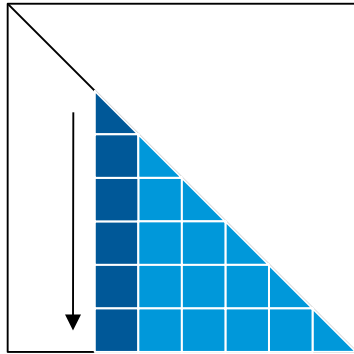


ICL

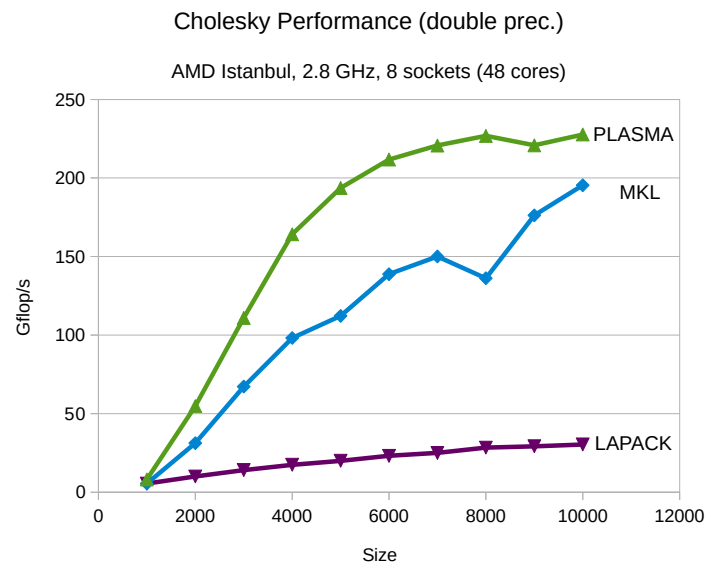
Algorithms

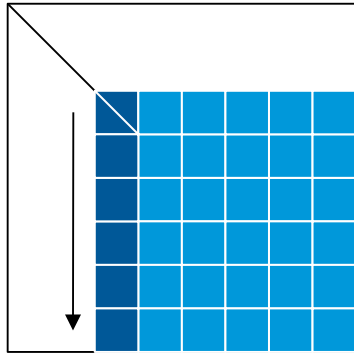
Cholesky

```
PLASMA_[scdz]potrf[_Tile][_Async]()
```



- **Algorithm**
 - equivalent to LAPACK
- **Numerics**
 - same as LAPACK
- **Performance**
 - comparable to vendor on few cores
 - much better than vendor on many cores



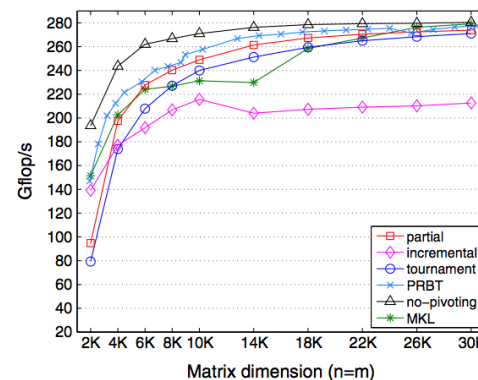


- **Algorithm**
 - equivalent to LAPACK
 - same pivot vector
 - same L and U factors
 - same forward substitution procedure

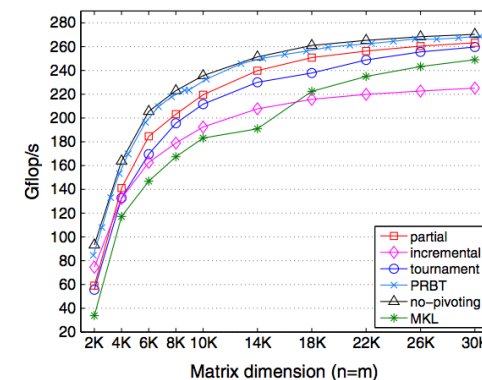
- **Numerics**
 - same as LAPACK

- **Performance**
 - comparable to vendor on few cores
 - much better than vendor on many cores

16 Sandy Bridge cores



Factorization alone, using 16 cores



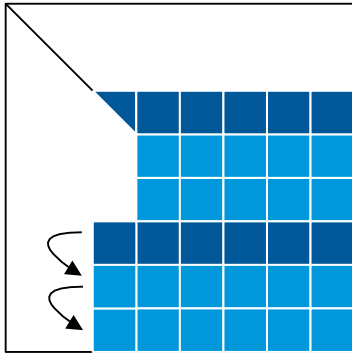
Factorization and solve with iterative refinement, using 16 cores



Algorithms

incremental QR Factorization

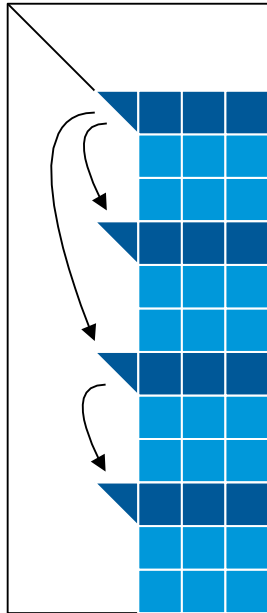
PLASMA_[scdz]geqrt[_Tile][_Async]()



- **Algorithm**
 - the same R factor as LAPACK (absolute values)
 - different set of Householder reflectors
 - different Q matrix
 - different Q generation / application procedure
- **Numerics**
 - same as LAPACK
- **Performance**
 - comparable to vendor on few cores
 - much better than vendor on many cores

```
PLASMA_[scdz]geqrt[_Tile][_Async]()
```

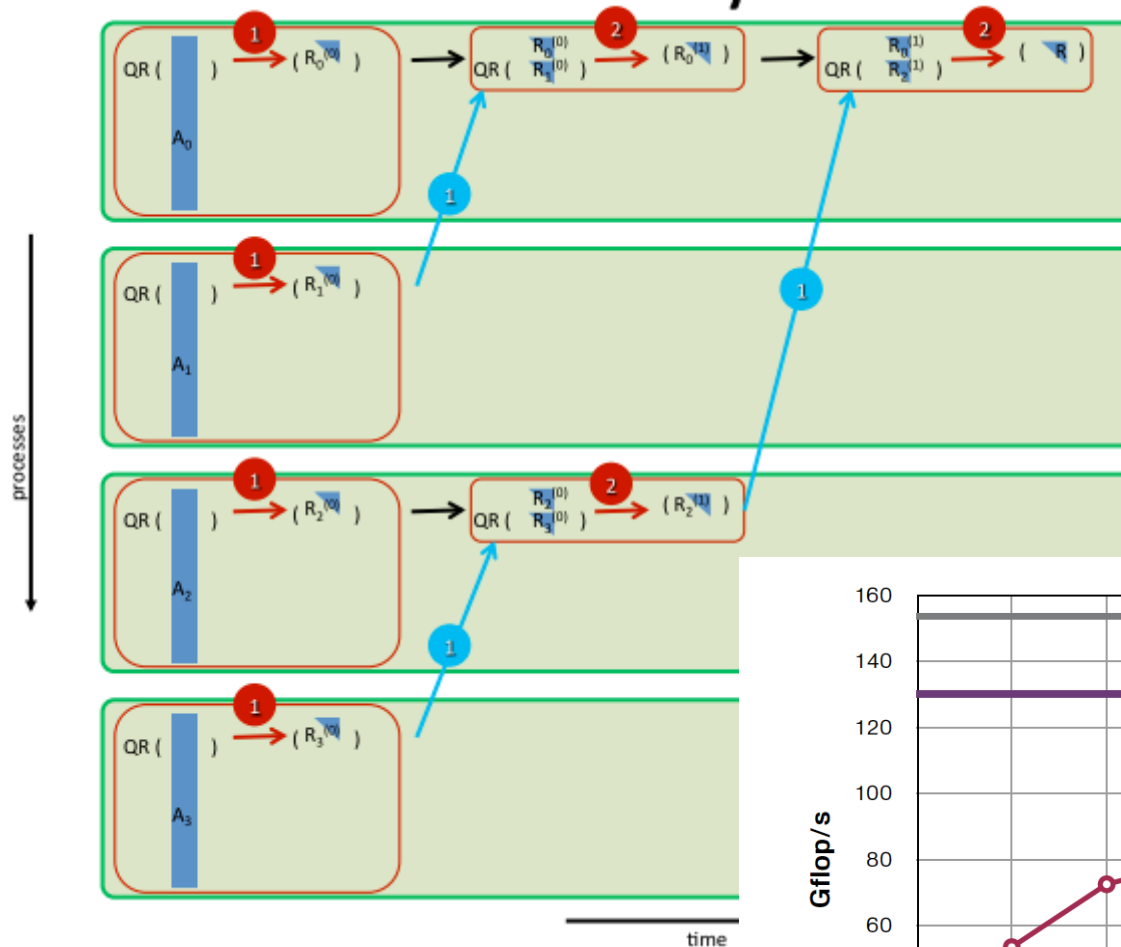
```
PLASMA_Set(
    PLASMA_HOUSEHOLDER_MODE,
    PLASMA_TREE_HOUSEHOLDER);
```



- **Algorithm**
 - the same R factor as LAPACK (absolute values)
 - different set of Householder reflectors
 - different Q matrix
 - different Q generation / application procedure
- **Numerics**
 - same as LAPACK
- **Performance**
 - absolutely superior for tall matrices

Communication Avoiding QR

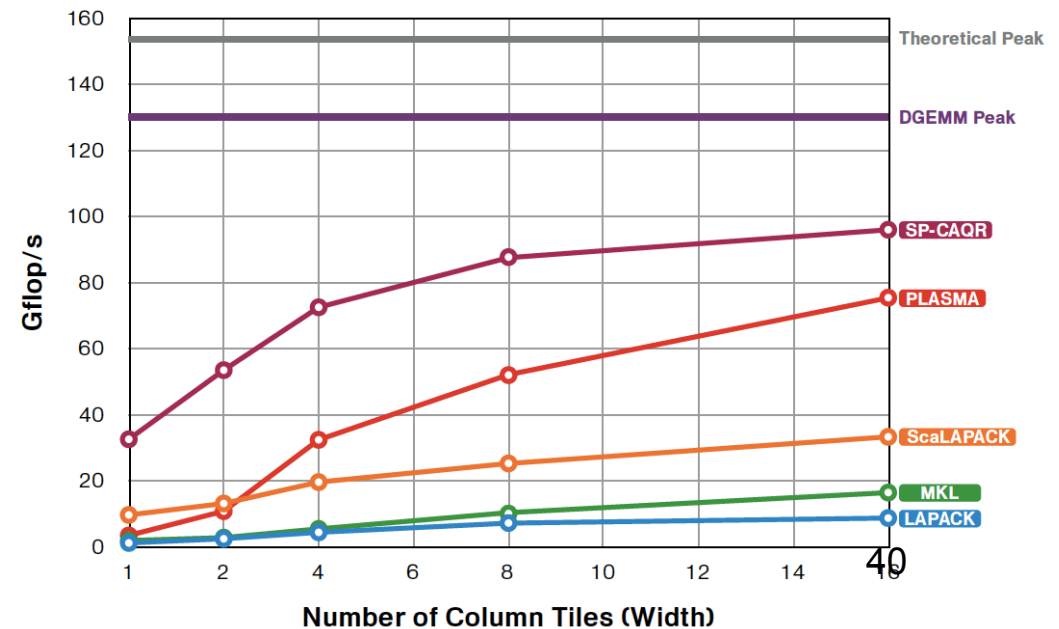
Example



$$Q_1^T \rightarrow Q_2^T \rightarrow Q_3^T \rightarrow R$$

$$A = Q_1 Q_2 Q_3 R = QR$$

Quad-socket, quad-core machine Intel Xeon EMT64 E7340 at 2.39 GHz.
Theoretical peak is 153.2 Gflop/s with 16 cores.
Matrix size 51200 by 3200



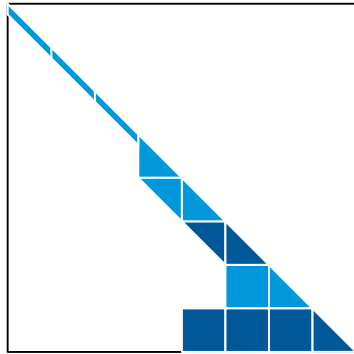


ICL

Algorithms

three-stage symmetric EVP

PLASMA_`[scdz]syev[_Tile][_Async]()`



- **Algorithm**

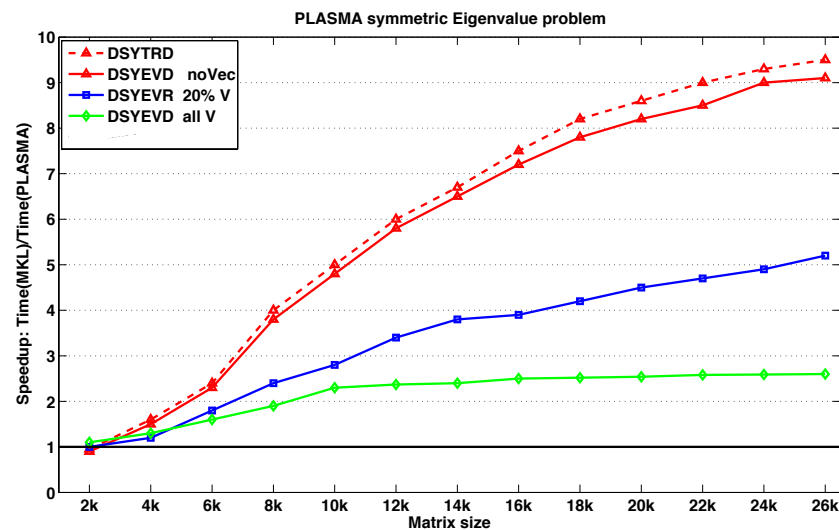
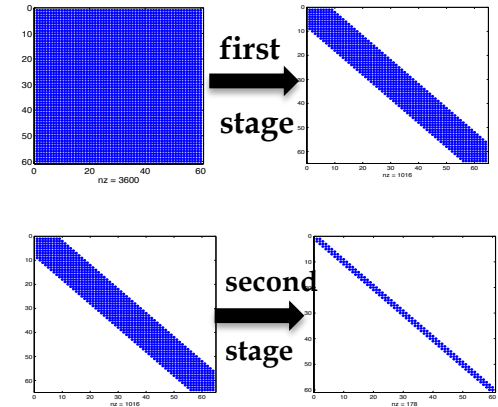
- two-stage tridiagonal reduction + QR Algorithm
- fast eigenvalues, slower eigenvectors
(possibility to calculate a subset)

- **Numerics**

- same as LAPACK

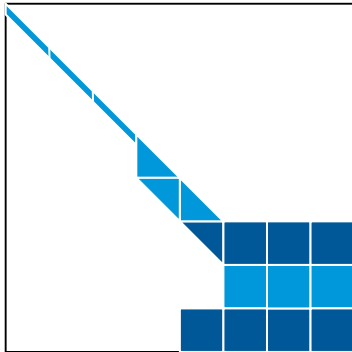
- **Performance**

- comparable to MKL for very small problems
- absolutely superior for larger problems



8/10/15

16 cores of Intel Sandy Bridge 41



Algorithm

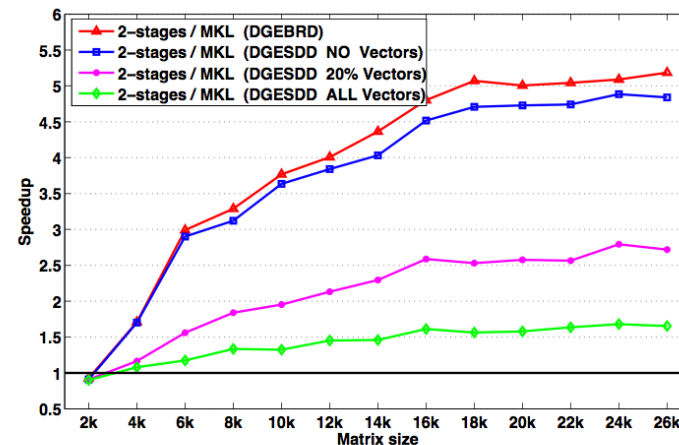
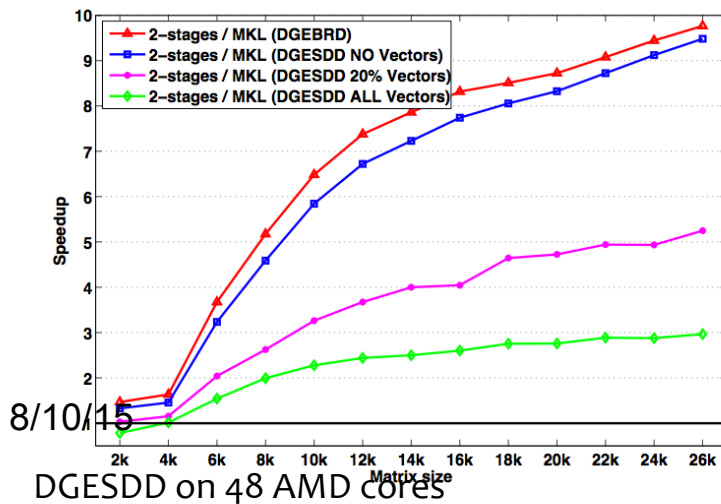
- two-stage bidiagonal reduction + QR iteration
- fast singular values, slower singular vectors
(possibility of calculating a subset)

Numerics

- same as LAPACK

Performance

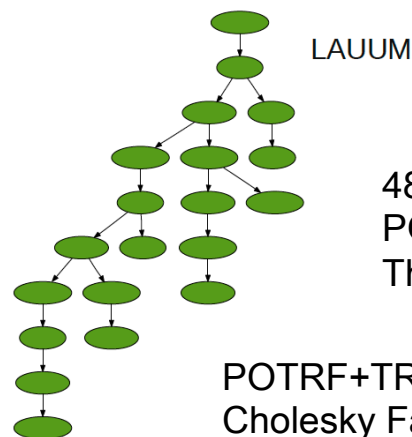
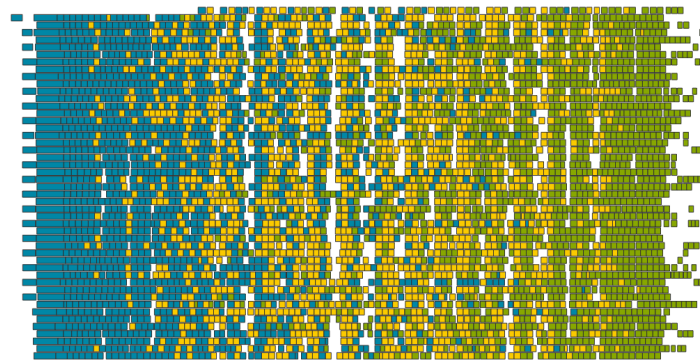
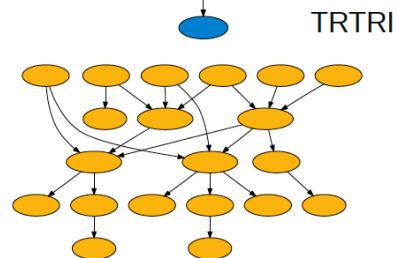
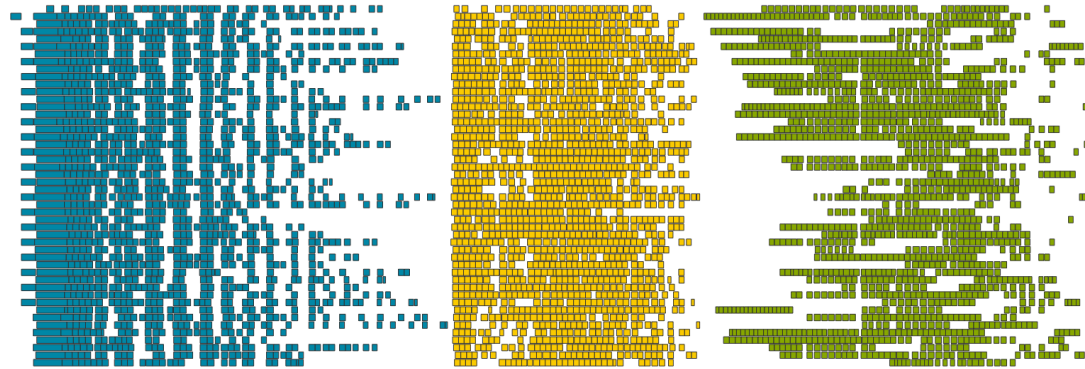
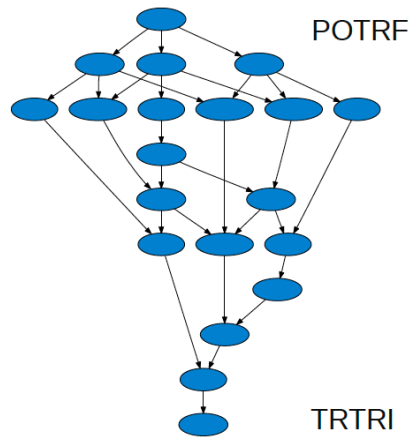
- comparable with MKL for very small problems
- absolutely superior for larger problems





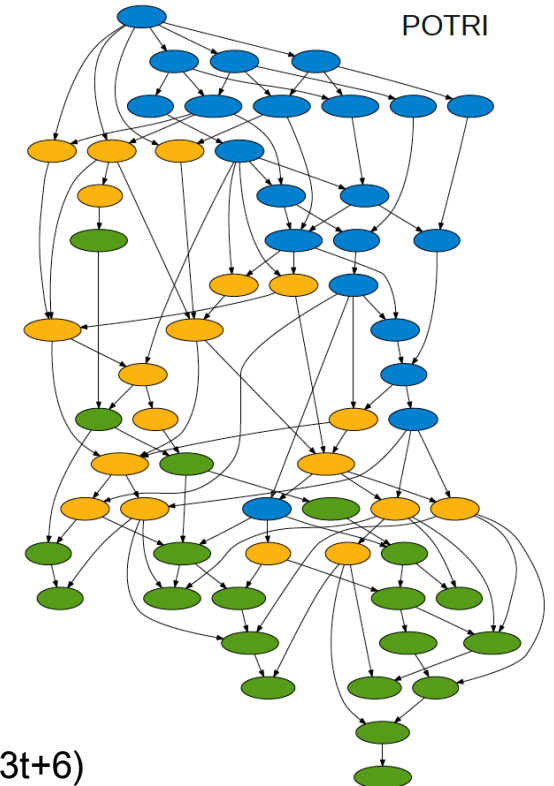
Pipelining: Cholesky Inversion

3 Steps: Factor, Invert L, Multiply L's



48 cores
POTRF, TRTRI and LAUUM.
The matrix is 4000 x 4000, tile size is 200 x 200,

POTRF+TRTRI+LAUUM: $25(7t-3)$
Cholesky Factorization alone: $3t-2$



Pipelined: $18(3t+6)$

Mixed Precision Methods

- **Mixed precision, use the lowest precision required to achieve a given accuracy outcome**
 - Improves runtime, reduce power consumption, lower data movement
 - Reformulate to find correction to solution, rather than solution; Δx rather than x .

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$
$$\boxed{x_{i+1} - x_i} = -\frac{f(x_i)}{f'(x_i)} \quad 44$$

Idea Goes Something Like This...

- Exploit 32 bit floating point as much as possible.
 - Especially for the bulk of the computation
- Correct or update the solution with selective use of 64 bit floating point to provide a refined results
- Intuitively:
 - Compute a 32 bit result,
 - Calculate a correction to 32 bit result using selected higher precision and,
 - Perform the update of the 32 bit results with the correction using high precision.

Mixed-Precision Iterative Refinement

- Iterative refinement for dense systems, $Ax = b$, can work this way.

$L U = \text{lu}(A)$	$O(n^3)$
$x = L \backslash (U \backslash b)$	$O(n^2)$
$r = b - Ax$	$O(n^2)$
WHILE $\ r\ $ not small enough	
$z = L \backslash (U \backslash r)$	$O(n^2)$
$x = x + z$	$O(n^1)$
$r = b - Ax$	$O(n^2)$
END	

- Wilkinson, Moler, Stewart, & Higham provide error bound for SP fl pt results when using DP fl pt.

Mixed-Precision Iterative Refinement

- Iterative refinement for dense systems, $Ax = b$, can work this way.

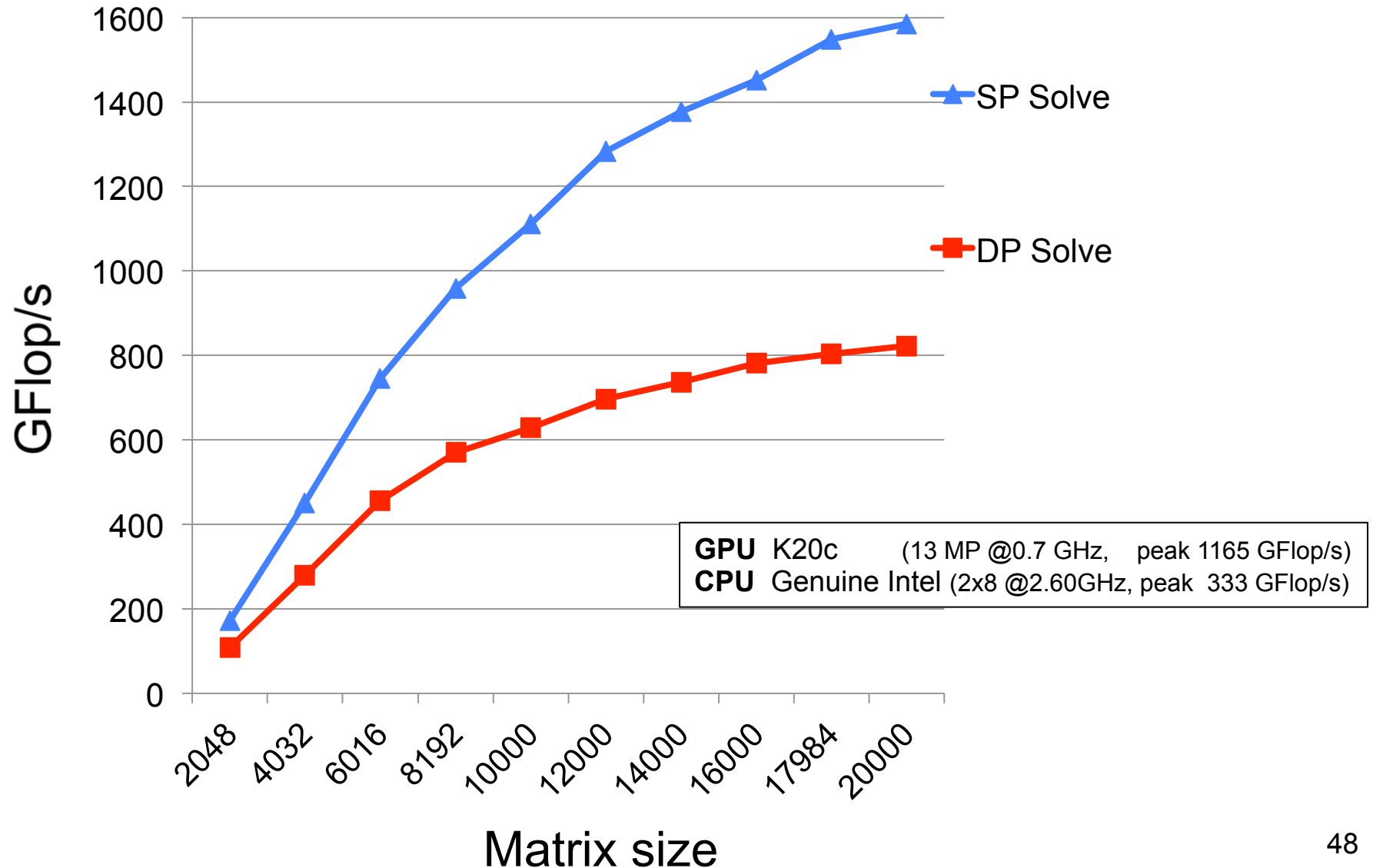
$L U = \text{lu}(A)$	SINGLE	$O(n^3)$
$x = L \backslash (U \backslash b)$	SINGLE	$O(n^2)$
$r = b - Ax$	DOUBLE	$O(n^2)$
WHILE $\ r\ $ not small enough		
$z = L \backslash (U \backslash r)$	SINGLE	$O(n^2)$
$x = x + z$	DOUBLE	$O(n^1)$
$r = b - Ax$	DOUBLE	$O(n^2)$
END		

- Wilkinson, Moler, Stewart, & Higham provide error bound for SP fl pt results when using DP fl pt.
- It can be shown that using this approach we can compute the solution to 64-bit floating point precision.

- Requires extra storage, total is 1.5 times normal;
- $O(n^3)$ work is done in lower precision
- $O(n^2)$ work is done in high precision
- Problems if the matrix is ill-conditioned in sp; $O(10^8)$

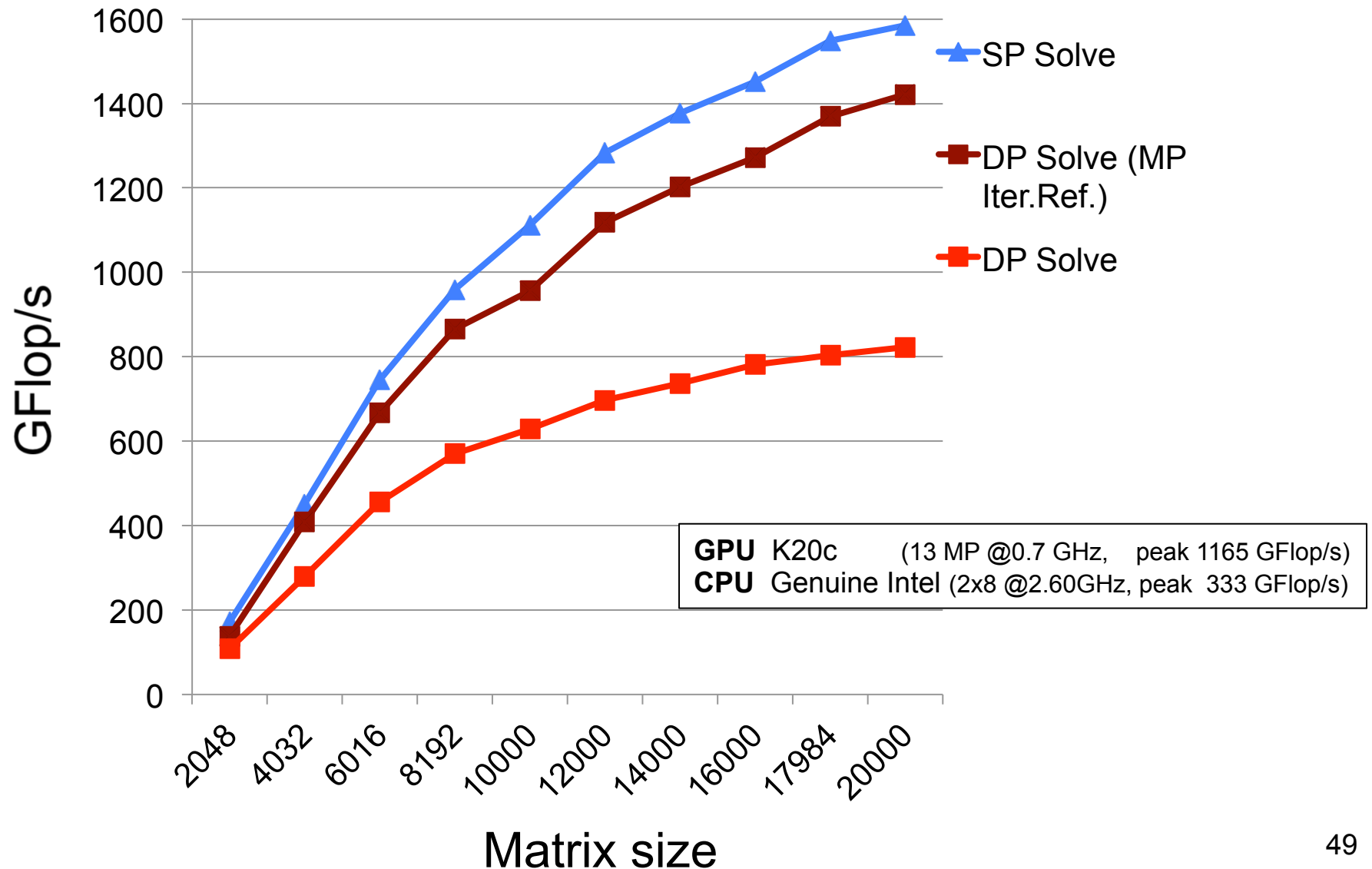
Mixed precision iterative refinement

Solving general dense linear systems using mixed precision iterative refinement



Mixed precision iterative refinement

Solving general dense linear systems using mixed precision iterative refinement





Critical Issues at Peta & Exascale for Algorithm and Software Design

- **Synchronization-reducing algorithms**
 - Break Fork-Join model
- **Communication-reducing algorithms**
 - Use methods which have lower bound on communication
- **Mixed precision methods**
 - 2x speed of ops and 2x speed for data movement
- **Autotuning**
 - Today's machines are too complicated, build “smarts” into software to adapt to the hardware
- **Fault resilient algorithms**
 - Implement algorithms that can recover from failures/bit flips
- **Reproducibility of results**
 - Today we can't guarantee this. We understand the issues, but some of our “colleagues” have a hard time with this.

Collaborators / Software / Support

- **PLASMA**
<http://icl.cs.utk.edu/plasma/>
- **MAGMA**
<http://icl.cs.utk.edu/magma/>
- **Quark (RT for Shared Memory)**
<http://icl.cs.utk.edu/quark/>
- **PaRSEC**(Parallel Runtime Scheduling
and Execution Control)
<http://icl.cs.utk.edu/parsec/>



- Collaborating partners
University of Tennessee, Knoxville
University of California, Berkeley
University of Colorado, Denver



